

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-073415

(43)Date of publication of application : 16.03.1999

(51)Int.Cl.

G06F 17/30

(21)Application number : 09-231363

(71)Applicant : TOSHIBA CORP

TOSHIBA COMPUT ENG CORP

(22)Date of filing : 27.08.1997

(72)Inventor : TANOSAKI YASUO

NAKAMOTO YUKIO

NISHINA TAKUYA

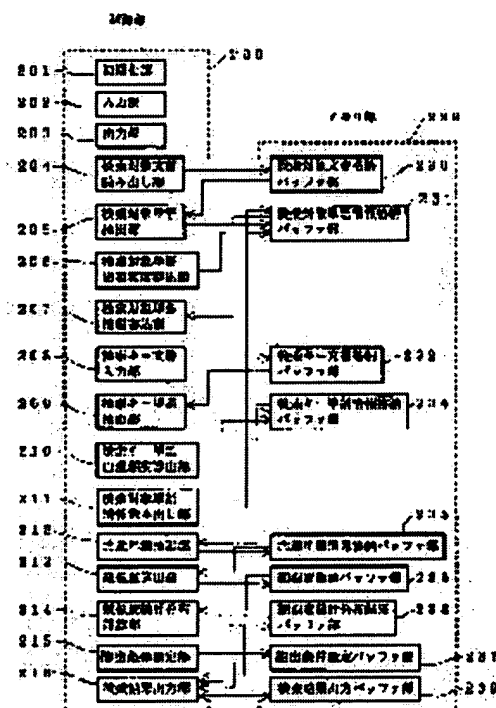
KUBOTA NAOHIDE

## (54) DEVICE AND METHOD FOR RETRIEVING SIMILAR DOCUMENT

### (57)Abstract:

**PROBLEM TO BE SOLVED:** To retrieve a similar document with high adequacy by finding statistical information on the similarity of each document to be retrieved and retrieving the similar document according to extraction conditions of the similar document set based on the similarity and statistical information.

**SOLUTION:** A common word extraction part 212 extracts the kind of a matching word and its appearance frequency information by comparing word information of a retrieval key document with word information of a document to be retrieved. A similarity calculation part 213 calculates the similarity between the retrieval key document and document to be retrieved according to the information. A similarity statistical distribution calculation part 214 finds statistical distribution calculation part 214 such as the mean value of similarity and standard deviation value. An extraction condition setting part 215 sets conditions for extracting the retrieval result of the similar document from the similarity statistical result. A retrieval result output part 216 judges whether or not there is the document to be retrieved and the corresponding document when there is the document to be retrieved from the



**BEST AVAILABLE COPY**

statistical distribution information, extraction condition values, and the similarity values of respective documents to be retrieved and outputs the retrieval result to a display device.

---

## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-73415

(43) 公開日 平成11年(1999) 3月16日

(51) Int.Cl.<sup>6</sup>  
G 0 6 F 17/30

識別記号

F I  
G 0 6 F 15/403

3 5 0 C

審査請求 未請求 請求項の数 4 O L (全 11 頁)

(21) 出願番号 特願平9-231363

(22) 出願日 平成9年(1997) 8月27日

(71) 出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(71) 出願人 000221052

東芝コンピュータエンジニアリング株式会  
社

東京都青梅市新町3丁目3番地の1

(72) 発明者 田野崎 康雄

東京都青梅市末広町2丁目9番地 株式会  
社東芝青梅工場内

(74) 代理人 弁理士 須山 佐一

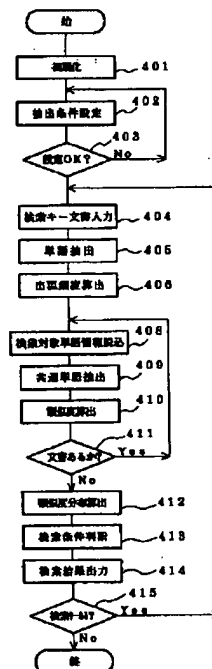
最終頁に続く

(54) 【発明の名称】 類似文書検索装置及び類似文書検索方法

(57) 【要約】

【課題】 ある文書（検索キー文書）と類似する文書を複数の検索対象文書のなかから検索する装置において、より信憑性の高い類似文書検索を実現する。

【解決手段】 検索キー文書と各検索対象文書との各々の類似度値の統計分布（例えば類似度の平均値）を求め、この統計分布を基準に、ユーザが設定した条件を満足するものを類似文書として抽出する。従来のように単に類似度値が高いものを類似文書として抽出する方式に比べ、類似文書としてより信憑性の高いものを検索結果として得ることができる。また、検索キー文書と各検索対象文書との類似度がどれも一般的な評価基準において高いとは言えないような場合に、類似文書がないことを検索結果として出力する。



【特許請求の範囲】

【請求項1】 検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索装置において、

前記検索キー文書と前記各検索対象文書との類似度を算出する類似度算出手段と、

前記類似度算出手段によって算出された各検索対象文書の類似度の統計情報を求める統計情報算出手段と、

前記統計情報を基準とする類似文書の抽出条件を設定する抽出条件設定手段と、

前記類似度算出手段によって算出された各検索対象文書の類似度および前記抽出条件設定手段により設定された抽出条件に基づいて類似文書を検索する検索手段とを具備することを特徴とする類似文書検索装置。

【請求項2】 検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索装置において、

前記検索キー文書と前記各検索対象文書との類似度を算出する類似度算出手段と、

前記類似度算出手段によって算出された各検索対象文書の類似度の統計情報を求める統計情報算出手段と、

前記統計情報を基準とする類似文書の有無の判定条件を設定する判定条件設定手段と、

前記類似度算出手段によって算出された各検索対象文書の類似度および前記判定条件設定手段により設定された判定条件に基づいて類似文書の有無を判定する判定手段とを具備することを特徴とする類似文書検索装置。

【請求項3】 検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索方法において、

前記検索キー文書と前記各検索対象文書との類似度を算出する工程と、

前記算出された各検索対象文書の類似度の統計情報を求める工程と、

前記統計情報を基準とする類似文書の抽出条件を設定する工程と、

前記算出された各検索対象文書の類似度および前記設定された抽出条件に基づいて類似文書を検索する工程とを具備することを特徴とする類似文書検索方法。

【請求項4】 検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索方法において、

前記検索キー文書と前記各検索対象文書との類似度を算出する工程と、

前記算出された各検索対象文書の類似度の統計情報を求める工程と、

前記統計情報を基準とする類似文書の有無の判定条件を設定する工程と、

前記算出された各検索対象文書の類似度および前記設定された判定条件に基づいて類似文書の有無を判定する工

程とを具備することを特徴とする類似文書検索方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、電子化された文書データの検索装置に係り、特にある文書データを検索キーとしてこれと類似した文書データを自動検索する類似文書検索装置および類似文書検索方法に関する。

【0002】

【従来の技術】近年、大量の電子化された文書データが流通するようになり、自動分類等を行う目的で、文書データベース中から指定された文書（以下、検索キー文書と呼ぶ。）に類似する文書の自動検索を行うシステムが実用されてきている。従来の類似文書検索システムでは、検索キー文書に含まれている単語と他の文書（以下、検索対象文書と呼ぶ。）に含まれている単語とを比較し、共通する単語の種類や出現回数・場所などからベクトル空間法により類似度を算出し、最も類似度の高い検索対象文書を検索結果として出力したり、類似度の高い文書から順に出力していた。

【0003】

【発明が解決しようとする課題】従来の類似文書検索方式は、検索キー文書と文書データベース中の各検索対象文書との類似度を各々算出し、より類似度の高い文書を判定する、例えば最大類似度のものを検索結果として出力している。しかし、その検索結果は必ずしも妥当なものとは言えない。すなわち、上記従来の類似文書検索方式により得た類似度に基づく検索結果は、各検索対象文書に対して求めた各類似度の単純な大小比較により得た検索結果にすぎないため、例えば、検索キー文書と特に類似している検索対象文書が1つも存在しないような場合でも、その中で最も類似度の高い検索対象文書を類似文書として無条件に出力してしまう。

【0004】本発明はこのような課題を解決するためのもので、複数の検索対象文書のなかから検索キー文書との類似が際立っているものを確実に検索することのできる類似文書検索装置および類似文書検索方法の提供を目的としている。

【0005】また、本発明は、多くの一般的な類似評価の基準において類似していると呼べる類似文書のみを検索結果として得ることで、信頼性の向上を図ることのできる類似文書検索装置および類似文書検索方法の提供を目的としている。

【0006】

【課題を解決するための手段】上記目的を達成するために、本発明の類似文書検索装置は、請求項1に記載されるように、検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索装置において、前記検索キー文書と前記各検索対象文書との類似度を算出する類似度算出手段と、前記類似度算出手段によって算出された各検索対象文書の類似度の統計情報を求める

統計情報算出手段と、前記統計情報を基準とする類似文書の抽出条件を設定する抽出条件設定手段と、前記類似度算出手段によって算出された各検索対象文書の類似度および前記抽出条件設定手段により設定された抽出条件に基づいて類似文書を検索する検索手段とを具備することを特徴とする。

【0007】また、本発明の類似文書検索装置は、請求項2に記載されるように、検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索装置において、前記検索キー文書と前記各検索対象文書との類似度を算出する類似度算出手段と、前記類似度算出手段によって算出された各検索対象文書の類似度の統計情報を求める統計情報算出手段と、前記統計情報を基準とする類似文書の有無の判定条件を設定する判定条件設定手段と、前記類似度算出手段によって算出された各検索対象文書の類似度および前記判定条件設定手段により設定された判定条件に基づいて類似文書の有無を判定する判定手段とを具備することを特徴とする。

【0008】さらに、本発明の類似文書検索方法は、請求項3に記載されるように、検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索方法において、前記検索キー文書と前記各検索対象文書との類似度を算出する工程と、前記算出された各検索対象文書の類似度の統計情報を求める工程と、前記統計情報を基準とする類似文書の抽出条件を設定する工程と、前記算出された各検索対象文書の類似度および前記設定された抽出条件に基づいて類似文書を検索する工程とを具備することを特徴とする。

【0009】さらに、本発明の類似文書検索方法は、請求項4に記載されるように、検索キー文書に類似する文書を複数の検索対象文書のなかから検索する類似文書検索方法において、前記検索キー文書と前記各検索対象文書との類似度を算出する工程と、前記算出された各検索対象文書の類似度の統計情報を求める工程と、前記統計情報を基準とする類似文書の有無の判定条件を設定する工程と、前記算出された各検索対象文書の類似度および前記設定された判定条件に基づいて類似文書の有無を判定する工程とを具備することを特徴とする。

【0010】請求項1および請求項3の発明によれば、各検索対象文書の類似度の統計情報を求め、各検索対象文書の類似度と、統計情報を基準に設定された類似文書の抽出条件に基づいて類似文書を検索することで、類似文書としてより妥当性の高いもの、つまり類似度がその他の多くの検索対象文書に比べ際立って高い検索対象文書を類似文書として検索することができる。

【0011】請求項2および請求項4の発明によれば、各検索対象文書の類似度の統計情報を求め、各検索対象文書の類似度と、統計情報を基準に設定された類似文書の有無の判定条件に基づいて類似文書の有無を判定することで、検索キー文書と各検索対象文書との類似度が

ずれも多く一般的な評価基準において高いと言えないような場合に類似文書が存在しないとし、一般的な評価基準において類似していると言える類似文書だけを検索結果として得ることができる。

【0012】

【発明の実施の形態】以下、本発明の実施の形態を図面を参照して詳細に説明する。

【0013】図1は本発明に係る一実施形態の類似文書検索装置のハードウェア構成を示す図である。

【0014】同図に示すように、この類似文書検索装置は、CPUおよびメモリなどから構成される制御装置1、キーボードなどの入力装置2、類似文書の検索結果などを表示する表示装置3、および文書データや類似文書検索のための各文書の単語情報などを格納する外部記憶装置4から構成される。

【0015】図2に本類似文書検索装置における制御装置1の構成を示す。制御装置1は制御部200とメモリ部229からなる。

【0016】制御部200は、初期化部201、入力部202、出力部203、検索対象文書読み出し部204、検索対象単語抽出部205、検索対象単語出現頻度算出部206、検索対象単語情報書込部207、検索キー文書入力部208、検索キー単語抽出部209、検索キー単語出現頻度算出部210、検索対象単語情報読み出し部211、共通単語抽出部212、類似度算出部213、類似度統計分布計算部214、抽出条件設定部215、検索結果出力部216などから構成される。メモリ部229は、検索対象文書格納バッファ部230、検索対象単語情報格納バッファ部231、検索キー文書格納バッファ部232、検索キー単語情報格納バッファ部234、共通単語情報格納バッファ部235、類似度格納バッファ部236、抽出条件設定バッファ部237、類似度統計分布結果バッファ部238、検索結果出力バッファ部239などから構成される。

【0017】初期化部201は、上記各バッファ部の初期化を行う。入力部202は、ユーザによる入力装置2からの検索キー文書や抽出条件の設定など各種設定の入力を行う。出力部203は、入力部202により入力された検索キー文書などの各種設定内容を表示装置3に出力する。

【0018】検索対象文書読み出し部204は、外部記憶装置4に格納されている検索対象文書に関する情報を文書データベース化するために、文書データベース化すべき文書情報を外部記憶装置4から読み込み、検索対象文書格納バッファ部230に格納する。

【0019】検索対象単語抽出部205は、検索対象文書格納バッファ部230に格納されている検索対象文書からの単語の切り出しを行う。そして、切り出した単語のなかからその文書内容を表す上でキーとなる単語を抽出し、抽出した単語種を検索対象単語情報格納バッファ

部231に格納する。単語の切り出しは形態素解析などにより行い、その文書の内容を表す上でキーとなる単語の単語種は品詞情報（例えば「名詞」や「サ変名詞」）を使って表現する。

【0020】検索対象単語出現頻度算出部206は、検索対象単語抽出部205により抽出された個々のキー単語について、検索対象文書中での出現頻度を算出し、これを検索対象文書の単語情報として検索対象単語情報格納バッファ部231に格納する。

【0021】検索対象単語情報書込部207は、検索対象単語情報格納バッファ部231に格納されている検索対象文書の単語情報を外部記憶装置4に格納する。

【0022】検索キー文書入力部208は、入力装置2から入力された検索キー文書の情報を検索キー文書格納バッファ部232に格納する。

【0023】検索キー単語抽出部209は、検索キー文書格納バッファ部232に格納されている検索キー文書からの単語切り出しを行う。そして、その文書の内容を表す上でキーとなる単語を抽出し、抽出した単語種を検索キー単語情報格納バッファ部234に格納する。単語の切り出しは形態素解析などにより行い、その文書の内容を表す上でキーとなる単語の単語種は品詞情報（例えば「名詞」や「サ変名詞」）を使って表現する。

【0024】検索キー単語出現頻度算出部210は、検索キー単語抽出部209により抽出された個々のキー単語について、検索キー文書中での出現頻度を算出し、これを検索キー文書の単語情報として検索キー単語情報格納バッファ部234に格納する。

【0025】検索対象単語情報読み出し部211は、外部記憶装置4に格納されている各検索対象文書の単語情報（単語の出現頻度情報）を1文書分ごとに呼び出し、検索対象単語情報格納バッファ部231に格納する。

【0026】共通単語抽出部212は、検索キー単語情報格納バッファ部234に格納されている検索キー文書の単語情報と検索対象単語情報格納バッファ部231に格納されている検索対象文書の単語情報とを比較して、一致する単語の種類と出現頻度情報を共通単語情報格納バッファ部235に格納する。

【0027】類似度算出部213は、共通単語情報格納バッファ部235に格納されている情報に基づき検索キー文書と検索対象文書との類似度を算出し、その類似度値を類似度格納バッファ部236に格納する。

【0028】類似度統計分布計算部214は、類似度格納バッファ部236に格納されている検索キー文書と全検索対象文書との類似度値から類似度の平均値や標準偏差値などの統計分布情報を求めて類似度統計分布結果バッファ部238に格納する。抽出条件設定部215は、入力装置2を介してユーザより入力された、類似度統計結果から類似文書の検索結果を抽出する場合の条件、または類似度統計結果から検索キー文書との類似文書があ

るとするための条件、または検索キー文書との類似文書がないとするための条件などの抽出条件値を抽出条件設定バッファ部237に格納（設定）する。

【0029】検索結果出力部216は、類似度統計分布結果バッファ部238に格納されている統計分布情報、抽出条件設定バッファ部237に格納されている抽出条件値、さらには類似度格納バッファ部236に格納されている各検索対象文書の類似度値から、検索キー文書に対する類似文書検索結果として、検索対象文書の有無、検索対象文書が有る場合の該当文書を判断し、その検索結果を検索結果出力バッファ部239に格納し、そして検索結果出力バッファ部239の内容を表示装置3に出力する。

【0030】次に、本実施形態の類似文書検索装置の動作を説明する。

【0031】最初に、検索対象文書のデータベースの作成手順を図3、図5、図6により説明する。図3はその手順を示すフローチャートである。

【0032】まず、初期化部201により全バッファ部の初期化を行う（ステップ301）。続いて検索対象文書読み出し部204が、外部記憶装置4から複数のテキスト文書を読み出し、検索対象文書格納バッファ部230に検索対象文書として格納する（ステップ302）。具体例として、例えば図5に示すような内容のテキスト文書を検索対象文書の一つとして格納したとする。

【0033】次に、検索対象単語抽出部205が、検索対象文書格納バッファ部230に格納されている個々の検索対象文書について、形態素解析などによって単語の切り出しを行い、切り出した単語のなかから文書内容を表すキー単語を抽出し、そのキー単語の単語種（例えば品詞情報）を検索対象単語情報格納バッファ部231に格納する（ステップ303）。

【0034】次に、検索対象単語出現頻度算出部206が、検索対象単語情報格納バッファ部231に格納されている検索対象文書のキー単語について、検索対象文書全体での出現頻度を算出し、その結果を検索対象単語情報格納バッファ部231に格納する（ステップ304）。図6に検索対象単語情報格納バッファ部231の格納例を示す。このバッファ部231において単語と頻度は対応付けて記述される。例えばキー単語「文書」が文書全体のなかで2回出現している場合は頻度として「2」が記述される。

【0035】このようにして検索対象単語情報格納バッファ部231に格納された情報は、検索対象文書のデータベースとして外部記憶装置4に蓄積される（ステップ305）。

【0036】この後、検索対象文書格納バッファ部230に文書データベース化前の検索対象文書が残っているかどうかを判断し（ステップ306）、他に検索対象文書があればステップ302に戻って、その新たな検索対

象文書についての前記同様の文書データベースの作成が行われる。他に検索対象文書がなければ本処理を終了する。

【0037】次に、類似文書の検索手順を図4、図7乃至図16により説明する。図4は類似文書検索手順を示すフローチャートである。

【0038】まず、初期化部201により全バッファ部の初期化を行う(ステップ401)。続いて抽出条件設定部215が起動される。抽出条件設定部215は、入力装置2を通じてユーザより、

1. 類似度統計結果から類似文書の検索結果を抽出する場合の条件
2. 類似度統計結果から類似文書が存在しているとするための条件
3. 類似度統計結果から類似文書が存在していないとするための条件

などの抽出条件値の入力を受け付けて抽出条件設定バッファ部237に格納(設定)する(ステップ402)。より具体的には、図7に抽出条件設定バッファ部237の格納例を示しているように、1.の条件として、「平均類似度の2倍以上の類似度を持つ検索対象文書を検索結果とする。」などが設定される。

【0039】2.の条件として、「平均類似度の2倍以上の類似度を持つ検索対象文書がある場合」などが設定される。

【0040】3.の条件として、「すべての検索対象文書の類似度が0.1以下である場合」などが設定される。

【0041】これらの抽出条件値はユーザにより任意に決定される。

【0042】この抽出条件の設定が完了したら(ステップ403)、検索キー文書入力部208が起動される。検索キー文書入力部208は、入力装置2を通じてユーザより検索キー文書の入力を受け付け、入力した検索キー文書の情報を検索キー文書格納バッファ部232に格納する(ステップ404)。具体例として、図8に示すような検索キー文書が入力されてバッファ部232に格納されたとする。

【0043】次に、検索キー単語抽出部209が、検索キー文書格納バッファ部232に格納されている検索キー文書から形態素解析などによって単語の切り出しを行い、切り出した単語のなかから文書内容を表すキーワードを抽出し、そのキーワードの単語種(例えば品詞情報)を検索キー単語情報格納バッファ部234に格納する(ステップ405)。

【0044】次に、検索キー単語出現頻度算出部210が、検索キー単語情報格納バッファ部234に格納されている個々のキーワードについて、検索キー文書全体の中での出現頻度を算出し、これを検索キー単語情報格納バッファ部234にキーワードと対応付けて格納する(ス

テップ406)。図9に検索キー単語情報格納バッファ部234の格納例を示す。このバッファ部234においてキーワードと頻度は対応付けて記述され、例えばキーワード「今後」が2回出現している場合は頻度として「2」が記述される。

【0045】次に、検索対象単語情報読み出し部211が、外部記憶装置4に格納されている各検索対象文書の単語情報を1文書ごとに読み込み、検索対象単語情報格納バッファ部231に書き込む(ステップ408)。

10 【0046】この後、共通単語抽出部212が起動され、共通単語抽出部212は、検索対象単語情報格納バッファ部231と検索キー単語情報格納バッファ部234とに共通に格納されているキーワードを検出し、共通単語情報格納バッファ部235に格納する(ステップ409)。例えば、図10に示すように、図6の検索対象単語情報格納バッファ部231と図9の検索キー単語情報格納バッファ部234に共通するキーワードとして「画像」が検出され、このキーワード「画像」とその頻度情報「3」を共通単語情報格納バッファ部235に対応付けて格納する。

20 【0047】次に、類似度算出部213が、共通単語情報格納バッファ部235に格納されている情報に基づき検索キー文書と検索対象文書との類似度をベクトル空間法などにより算出し、その類似度値を類似度格納バッファ部236に格納する(ステップ410)。例えば、図11に示すような各検索対象文書ごとの類似度が類似度格納バッファ部236に格納される。

【0048】すべての検索対象文書について類似度計算が完了すると(ステップ411)、類似度統計分布計算部214が起動する。類似度統計分布計算部214は、類似度格納バッファ部236に格納されている各検索対象文書の類似度の統計分布を算出し、その結果を類似度統計分布結果バッファ部238に格納する(ステップ412)。例えば、図12に示すように、各検索対象文書の類似度の平均値「0.25」を類似度の統計分布情報として求める。

【0049】次に、検索結果出力部216が、抽出条件設定バッファ部237に格納されている抽出条件値を判断し(ステップ413)、その抽出条件値、類似度統計分布結果バッファ部238に格納されている統計分布情報、さらには類似度格納バッファ部236に格納されている各検索対象文書の類似度値から、検索キー文書に対する類似文書の検索結果として、検索対象文書の有無や、検索対象文書が有る場合の該当文書を判断し、その結果を検索結果出力バッファ部239に格納する。例えば、1.の条件(類似度統計結果から類似文書の検索結果を抽出する場合の条件)に対し、図13に示すように、上記条件を満足する検索対象文書があれば、そのIDを検索結果出力バッファ部239に格納する。

50 【0050】また、2.の条件(類似度統計結果から検

索キー文書との類似文書が存在しているとするための条件) に対し、該条件を満足している場合は、図 14 に示すように、当該検索キー文書を類似文書有りの検索キー文書として、その ID を検索結果出力バッファ部 239 に格納する。

【0051】さらに、3. の条件 (類似度統計結果から検索キー文書との類似文書が存在していないとするための条件) に対し、該条件を満足する場合は、図 15 に示すように、当該検索キー文書を類似文書無しの検索キー文書として、その ID を検索結果出力バッファ部 239 に格納する。

【0052】検索結果出力部 216 は、検索結果出力バッファ部 239 の内容を、例えば図 16 に示すような形式で表示装置 3 に出力する (ステップ 414)。図 16 の例では、例えば図 14 に示す ID「2」の類似文書有りの検索キー文書と類似する検索対象文書として ID「1」「42」「54」「314」などがあることを示している。また、2. と 3. の各条件により、類似する検索対象文書がないと判断された検索キー文書に対しては類似文書がないことが表示される。

【0053】この後、次の検索キー文書がある場合はステップ 404 に戻ってその検索キー文書を入力し、以降同様の処理を行う。次の検索キー文書がなければ本処理を終了する。

【0054】なお、本動作例では、1つの検索キー文書に対する検索結果を得たところでこれを表示するようにしたが、すべての検索キー文書に対する検索結果を得た後、検索結果を見たい検索キー文書を指定すると、その検索結果が表示されるようにしてもよい。

【0055】このように本実施形態の類似文書検索装置においては、検索キー文書と各検索対象文書との各々の類似度値の統計分布 (この実施形態では類似度の平均値) を求め、この統計分布を基準に、ユーザが設定した条件を満足するものを類似文書として抽出することで、従来のように単に類似度値が高いものを類似文書として抽出する方式に比べ、類似文書としてより信憑性の高いものを検索結果として得ることができる。すなわち、本実施形態は、類似度の統計分布を基準としているので、検索キー文書との類似度がその他多くの検索対象文書に比べ際立って高い検索対象文書を類似文書として得られる。

【0056】また、従来の方式では、検索キー文書と各検索対象文書との類似度が、どれも一般的な評価基準において高いとは言えないような場合でも検索結果として類似文書は無条件に出力してしまうが、本実施形態では、このような場合において類似文書がないことを検索結果として出力する。このような点からも、本実施形態の類似文書検索装置によれば、類似文書として信憑性の高い検索結果を得ることができ、類似文書検索効率を大幅に高めることができる。

【0057】

【発明の効果】以上説明したように本発明によれば、各検索対象文書の類似度の統計情報を求め、各検索対象文書の類似度と、統計情報を基準に設定された類似文書の抽出条件に基づいて類似文書を検索することで、類似文書としてより妥当性の高いもの、つまり類似度がその他の多くの検索対象文書に比べ際立って高い検索対象文書を類似文書として得ることができる。

【0058】また、本発明によれば、各検索対象文書の類似度の統計情報を求め、各検索対象文書の類似度と、統計情報を基準に設定された類似文書の有無の判定条件に基づいて類似文書の有無を判定することで、検索キー文書と各検索対象文書との類似度がいずれも一般的な評価基準において高いと言えないような場合に類似文書が存在しないとし、一般的な評価基準において類似していると言える類似文書だけを検索結果として得ることができる。

【図面の簡単な説明】

【図 1】本発明に係る一実施形態の類似文書検索装置のハードウェア構成を示す図

【図 2】図 1 の類似文書検索装置における制御装置の機能ブロック図

【図 3】検索対象文書のデータベースの作成手順を示す図

【図 4】類似文書検索手順を示す図

【図 5】検索対象文書の例

【図 6】検索対象単語情報の格納例

【図 7】抽出条件の設定例

【図 8】検索キー文書の例

【図 9】検索キー単語情報の例

【図 10】共通単語情報の例

【図 11】各検索対象文書の類似度の例

【図 12】類似度の平均値の例

【図 13】類似文書の検索結果の例

【図 14】類似文書の検索結果の例

【図 15】類似文書の検索結果の例

【図 16】類似文書の検索結果の出力例

【符号の説明】

200……制御部

210……初期化部

202……入力部

203……出力部

204……検索対象文書読み出し部

205……検索対象単語抽出部

206……検索対象単語出現頻度算出部

207……検索対象単語情報書込部

208……検索キー文書入力部

209……検索キー単語抽出部

210……検索キー単語出現頻度算出部

50 211……検索対象単語情報読み出し部



11

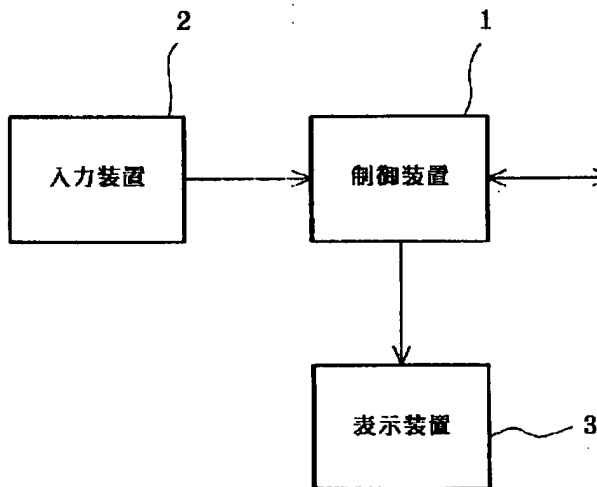
212……共通単語抽出部  
 213……類似度算出部  
 214……類似度統計分布計算部  
 215……抽出条件設定部  
 216……検索結果出力部  
 229……メモリ部  
 230……検索対象文書格納バッファ部  
 231……検索対象単語情報格納バッファ部

\*

12

\*232……検索キー文書格納バッファ部  
 233……検索キー単語情報格納バッファ部  
 235……共通単語情報格納バッファ部  
 236……類似度格納バッファ部  
 237……抽出条件設定バッファ部  
 238……類似度統計分布結果バッファ部  
 239……検索結果出力バッファ部

【図1】



【図5】

この文書は、画像について書いてあります。  
 ……

検索対象文書バッファ格納例

【図8】

今後は、画像に関する事項が  
 ……

検索キー文書バッファ格納例

【図6】

単語	類似度
文書	2
画像	3
……	……

検索対象単語情報バッファ格納例

【図7】

抽出検索対象文書＝ 平均類似度の2倍以上
類似検索対象文書有り＝ 平均類似度の2倍以上の文書がある
類似度検索対象文書無し＝ すべての検索対象文書の類似度が0.1以下

抽出条件設定バッファ格納例

【図9】

単語	類似度
今後	2
画像	3
……	……

検索キー単語情報バッファ格納例

【図10】

単語	類似度
単語	3
……	……

共通単語情報バッファ格納例

【図12】

平均値＝0.25
----------

類似度統計分布結果バッファ格納例

【図14】

【図11】

検索対象文書ID	類似度
1	0.2464
2	0.6842
3	0.9542
……	……

類似度バッファ格納例

【図13】

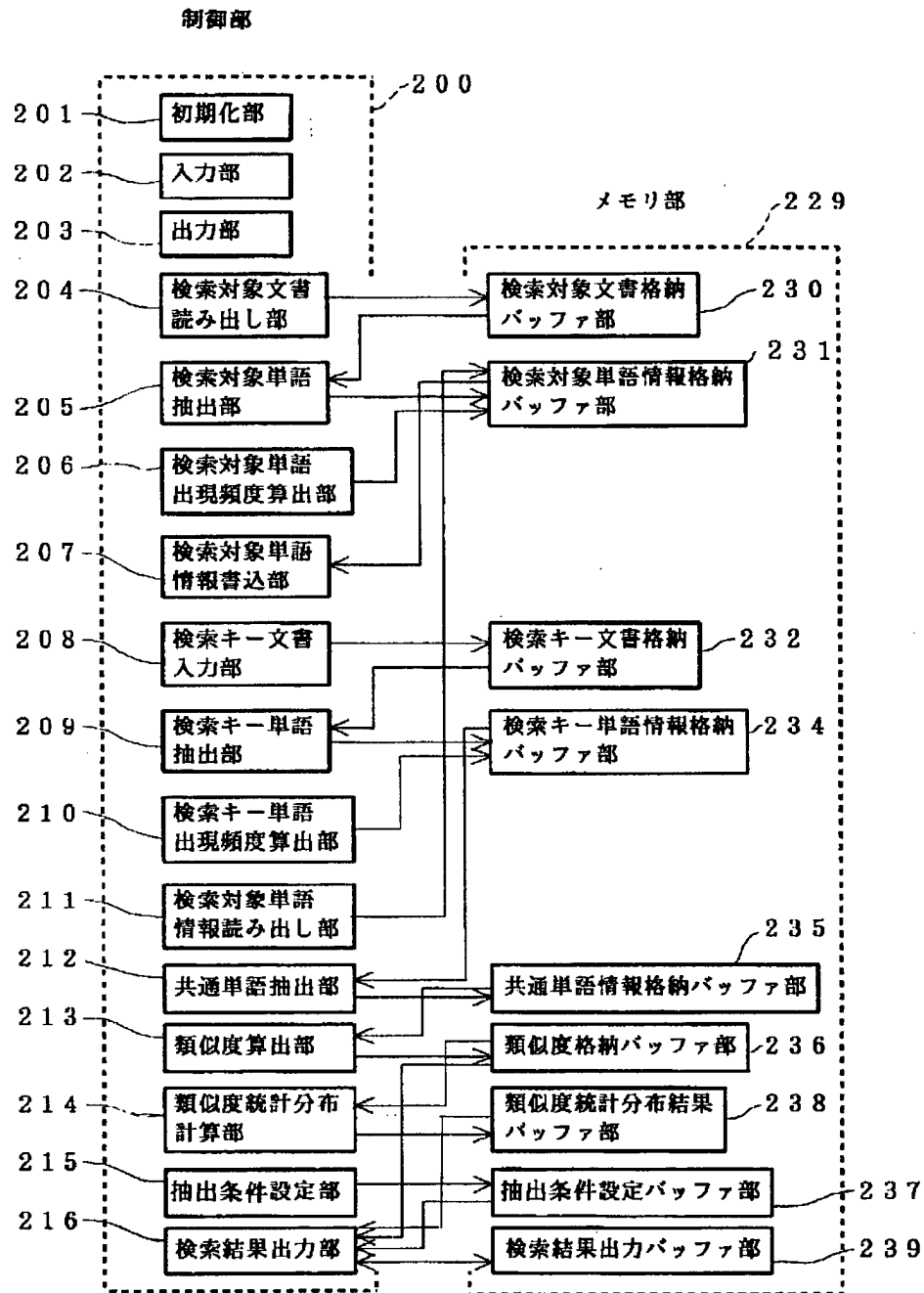
類似検索対象文書＝
1
42
54
314
……

検索結果出力バッファ格納例

類似文書有り検索キー文書＝
2
3
……

検索結果出力バッファ格納例

【図2】



【図15】

【図16】

類似文書無し検索キー文書＝

4

6

.....

検索結果出力バッファ格納例

結果：類似検索対象文書

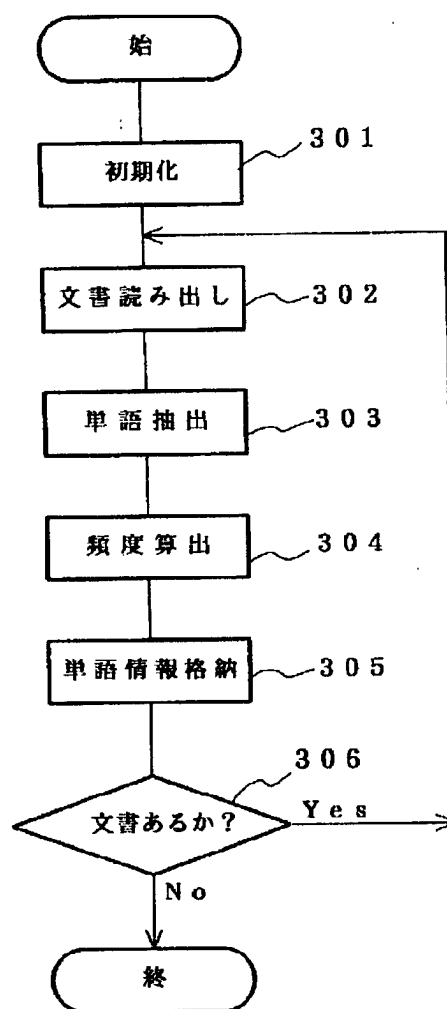
1

42

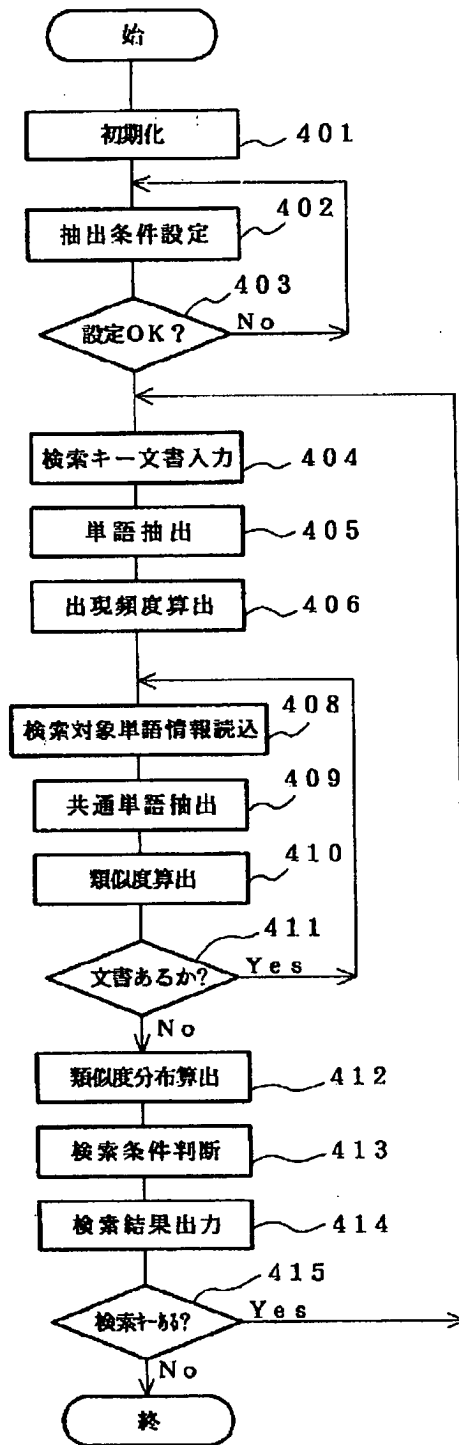
54

314

【図3】



【図4】



フロントページの続き

(72)発明者 中本 幸夫  
東京都青梅市新町1381番地 1 東芝コンピ  
ュータエンジニアリング株式会社内

(72)発明者 仁科 卓哉  
東京都青梅市新町1381番地 1 東芝コンピ  
ュータエンジニアリング株式会社内

(72)発明者 久保田 直秀  
東京都青梅市新町1381番地 1 東芝コンピ  
ュータエンジニアリング株式会社内

**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

☒ **BLACK BORDERS**

☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**

☐ **FADED TEXT OR DRAWING**

☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**

☐ **SKEWED/SLANTED IMAGES**

☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**

☐ **GRAY SCALE DOCUMENTS**

☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**

☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**

☐ **OTHER:** \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**